



## sklearn.tree.DecisionTreeClassifier

»

```
class sklearn.tree.DecisionTreeClassifier(criterion='gini', splitter='best', max_depth=None,
min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None,
random_state=None, max_leaf_nodes=None, class_weight=None, presort=False) \[source\]
```

A decision tree classifier.

Read more in the [User Guide](#).

**Parameters:** **criterion** : string, optional (default="gini")

The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

**splitter** : string, optional (default="best")

The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.

**max\_features** : int, float, string or None, optional (default=None)

The number of features to consider when looking for the best split:

- If int, then consider *max\_features* features at each split.
- If float, then *max\_features* is a percentage and  $\text{int}(\text{max\_features} * n\_features)$  features are considered at each split.
- If "auto", then  $\text{max\_features} = \text{sqrt}(n\_features)$ .
- If "sqrt", then  $\text{max\_features} = \text{sqrt}(n\_features)$ .
- If "log2", then  $\text{max\_features} = \log_2(n\_features)$ .
- If None, then  $\text{max\_features} = n\_features$ .

Note: the search for a split does not stop until at least one valid partition of the node samples is found, even if it requires to effectively inspect more than *max\_features* features.

**max\_depth** : int or None, optional (default=None)

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than *min\_samples\_split* samples. Ignored if *max\_leaf\_nodes* is not None.

**min\_samples\_split** : int, optional (default=2)

The minimum number of samples required to split an internal node.

**min\_samples\_leaf** : int, optional (default=1)

The minimum number of samples required to be at a leaf node.

**min\_weight\_fraction\_leaf** : float, optional (default=0.)

The minimum weighted fraction of the input samples required to be at a leaf node.

»

**max\_leaf\_nodes** : int or None, optional (default=None)

Grow a tree with max\_leaf\_nodes in best-first fashion. Best nodes are defined as relative reduction in impurity. If None then unlimited number of leaf nodes. If not None then max\_depth will be ignored.

**class\_weight** : dict, list of dicts, “balanced” or None, optional (default=None)

Weights associated with classes in the form {class\_label: weight}. If not given, all classes are supposed to have weight one. For multi-output problems, a list of dicts can be provided in the same order as the columns of y.

The “balanced” mode uses the values of y to automatically adjust weights inversely proportional to class frequencies in the input data as  $n_{\text{samples}} / (n_{\text{classes}} * \text{np.bincount}(y))$

For multi-output, the weights of each column of y will be multiplied.

Note that these weights will be multiplied with sample\_weight (passed through the fit method) if sample\_weight is specified.

**random\_state** : int, RandomState instance or None, optional (default=None)

If int, random\_state is the seed used by the random number generator; If RandomState instance, random\_state is the random number generator; If None, the random number generator is the RandomState instance used by *np.random*.

**presort** : bool, optional (default=False)

Whether to presort the data to speed up the finding of best splits in fitting. For the default settings of a decision tree on large datasets, setting this to true may slow down the training process. When using either a smaller dataset or a restricted depth, this may speed up the training.

---

**Attributes:**    **classes\_** : array of shape = [n\_classes] or a list of such arrays

The classes labels (single output problem), or a list of arrays of class labels (multi-output problem).

**feature\_importances\_** : array of shape = [n\_features]

The feature importances. The higher, the more important the feature. The importance of a feature is computed as the (normalized) total reduction of the criterion brought by that feature. It is also known as the Gini importance [R66].

**max\_features\_** : int,

The inferred value of max\_features.

»

**n\_classes\_** : int or list

The number of classes (for single output problems), or a list containing the number of classes for each output (for multi-output problems).

**n\_features\_** : int

The number of features when fit is performed.

**n\_outputs\_** : int

The number of outputs when fit is performed.

**tree\_** : Tree object

The underlying Tree object.

**See also:** [DecisionTreeRegressor](#)

## References

[R63] [http://en.wikipedia.org/wiki/Decision\\_tree\\_learning](http://en.wikipedia.org/wiki/Decision_tree_learning)

[R64] L. Breiman, J. Friedman, R. Olshen, and C. Stone, “Classification and Regression Trees”, Wadsworth, Belmont, CA, 1984.

[R65] T. Hastie, R. Tibshirani and J. Friedman. “Elements of Statistical Learning”, Springer, 2009.

[R66] (1, 2) L. Breiman, and A. Cutler, “Random Forests”, [http://www.stat.berkeley.edu/~breiman/RandomForests/cc\\_home.htm](http://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.htm)

## Examples

```
>>> from sklearn.datasets import load_iris
>>> from sklearn.cross_validation import cross_val_score
>>> from sklearn.tree import DecisionTreeClassifier
>>> clf = DecisionTreeClassifier(random_state=0)
>>> iris = load_iris()
>>> cross_val_score(clf, iris.data, iris.target, cv=10)
...
array([ 1.    ,  0.93...,  0.86...,  0.93...,  0.93...,
        0.93...,  0.93...,  1.    ,  0.93...,  1.    ])
```

&gt;&gt;&gt;

## Methods

**apply**(X[, check\_input])

Returns the index of the leaf that each sample is predicted as.

<code>fit(X, y[, sample_weight, check_input, ...])</code>	Build a decision tree from the training set (X, y).
<code>fit_transform(X[, y])</code>	Fit to data, then transform it.
<code>get_params([deep])</code>	Get parameters for this estimator.
<code>predict(X[, check_input])</code>	Predict class or regression value for X.
<code>predict_log_proba(X)</code>	Predict class log-probabilities of the input samples X.
<code>predict_proba(X[, check_input])</code>	Predict class probabilities of the input samples X.
<code>score(X, y[, sample_weight])</code>	Returns the mean accuracy on the given test data and labels.
<code>set_params(**params)</code>	Set the parameters of this estimator.
<code>transform(*args, **kwargs)</code>	DEPRECATED: Support to use estimators as feature selectors will be removed in version 0.19.

`__init__(criterion='gini', splitter='best', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, class_weight=None, presort=False)` [\[source\]](#)

`apply(X, check_input=True)` [\[source\]](#)

Returns the index of the leaf that each sample is predicted as.

*New in version 0.17.*

**Parameters:** **X** : array\_like or sparse matrix, shape = [n\_samples, n\_features]

The input samples. Internally, it will be converted to dtype=np.float32 and if a sparse matrix is provided to a sparse csr\_matrix.

**check\_input** : boolean, (default=True)

Allow to bypass several input checking. Don't use this parameter unless you know what you do.

**Returns:** **X\_leaves** : array\_like, shape = [n\_samples,]

For each datapoint x in X, return the index of the leaf x ends up in. Leaves are numbered within [0; self.tree\_.node\_count), possibly with gaps in the numbering.

## feature\_importances\_

Return the feature importances.

The importance of a feature is computed as the (normalized) total reduction of the criterion brought by that feature. It is also known as the Gini importance.

**Returns:** **feature\_importances\_** : array, shape = [n\_features]

**fit**(*X*, *y*, *sample\_weight=None*, *check\_input=True*, *X\_idx\_sorted=None*)

[\[source\]](#)

Build a decision tree from the training set (*X*, *y*).

**Parameters:** **X** : array-like or sparse matrix, shape = [*n\_samples*, *n\_features*]

The training input samples. Internally, it will be converted to dtype=np.float32 and if a sparse matrix is provided to a sparse csc\_matrix.

»

**y** : array-like, shape = [*n\_samples*] or [*n\_samples*, *n\_outputs*]

The target values (class labels in classification, real numbers in regression). In the regression case, use dtype=np.float64 and order='C' for maximum efficiency.

**sample\_weight** : array-like, shape = [*n\_samples*] or None

Sample weights. If None, then samples are equally weighted. Splits that would create child nodes with net zero or negative weight are ignored while searching for a split in each node. In the case of classification, splits are also ignored if they would result in any single class carrying a negative weight in either child node.

**check\_input** : boolean, (default=True)

Allow to bypass several input checking. Don't use this parameter unless you know what you do.

**X\_idx\_sorted** : array-like, shape = [*n\_samples*, *n\_features*], optional

The indexes of the sorted training input samples. If many tree are grown on the same dataset, this allows the ordering to be cached between trees. If None, the data will be sorted here. Don't use this parameter unless you know what to do.

**Returns:** **self** : object

Returns self.

**fit\_transform**(*X*, *y=None*, *\*\*fit\_params*)

[\[source\]](#)

Fit to data, then transform it.

Fits transformer to *X* and *y* with optional parameters *fit\_params* and returns a transformed version of *X*.

**Parameters:** **X** : numpy array of shape [*n\_samples*, *n\_features*]

Training set.

**y** : numpy array of shape [n\_samples]

Target values.

---

**Returns:**     **X\_new** : numpy array of shape [n\_samples, n\_features\_new]

Transformed array.

---

»

**get\_params(*deep=True*)**

[\[source\]](#)

Get parameters for this estimator.

---

**Parameters:**   **deep: boolean, optional :**

If True, will return the parameters for this estimator and contained subobjects that are estimators.

---

**Returns:**     **params** : mapping of string to any

Parameter names mapped to their values.

---

**predict(*X, check\_input=True*)**

[\[source\]](#)

Predict class or regression value for X.

For a classification model, the predicted class for each sample in X is returned. For a regression model, the predicted value based on X is returned.

---

**Parameters:**   **X** : array-like or sparse matrix of shape = [n\_samples, n\_features]

The input samples. Internally, it will be converted to dtype=np.float32 and if a sparse matrix is provided to a sparse csr\_matrix.

**check\_input** : boolean, (default=True)

Allow to bypass several input checking. Don't use this parameter unless you know what you do.

---

**Returns:**     **y** : array of shape = [n\_samples] or [n\_samples, n\_outputs]

The predicted classes, or the predict values.

---

**predict\_log\_proba(*X*)**

[\[source\]](#)

Predict class log-probabilities of the input samples X.

---

**Parameters:**   **X** : array-like or sparse matrix of shape = [n\_samples, n\_features]

The input samples. Internally, it will be converted to `dtype=np.float32` and if a sparse matrix is provided to a sparse `csr_matrix`.

---

**Returns:** **p** : array of shape = [n\_samples, n\_classes], or a list of n\_outputs

such arrays if `n_outputs > 1`. The class log-probabilities of the input samples. The order of the classes corresponds to that in the attribute `classes_`.

»

---

**predict\_proba**(*X*, *check\_input=True*)

[\[source\]](#)

Predict class probabilities of the input samples *X*.

The predicted class probability is the fraction of samples of the same class in a leaf.

`check_input` : *boolean*, (*default=True*)

Allow to bypass several input checking. Don't use this parameter unless you know what you do.

---

**Parameters:** **X** : array-like or sparse matrix of shape = [n\_samples, n\_features]

The input samples. Internally, it will be converted to `dtype=np.float32` and if a sparse matrix is provided to a sparse `csr_matrix`.

---

**Returns:** **p** : array of shape = [n\_samples, n\_classes], or a list of n\_outputs

such arrays if `n_outputs > 1`. The class probabilities of the input samples. The order of the classes corresponds to that in the attribute `classes_`.

---

**score**(*X*, *y*, *sample\_weight=None*)

[\[source\]](#)

Returns the mean accuracy on the given test data and labels.

In multi-label classification, this is the subset accuracy which is a harsh metric since you require for each sample that each label set be correctly predicted.

---

**Parameters:** **X** : array-like, shape = (n\_samples, n\_features)

Test samples.

**y** : array-like, shape = (n\_samples) or (n\_samples, n\_outputs)

True labels for *X*.

**sample\_weight** : array-like, shape = [n\_samples], optional

Sample weights.

---

**Returns:**     **score** : float

Mean accuracy of self.predict(X) wrt. y.

---

**set\_params**(\*\*params)

[\[source\]](#)

»

Set the parameters of this estimator.

The method works on simple estimators as well as on nested objects (such as pipelines). The former have parameters of the form <component>\_\_<parameter> so that it's possible to update each component of a nested object.

---

**Returns:**   **self** :

---

**transform**(\*args, \*\*kwargs)

[\[source\]](#)

DEPRECATED: Support to use estimators as feature selectors will be removed in version 0.19. Use SelectFromModel instead.

Reduce X to its most important features.

Uses coef\_ or feature\_importances\_ to determine the most important features. For models with a coef\_ for each class, the absolute sum over the classes is used.

---

**Parameters:**   **X** : array or scipy sparse matrix of shape [n\_samples, n\_features]

The input samples.

**threshold** : *string, float or None, optional (default=None)*

The threshold value to use for feature selection. Features whose importance is greater or equal are kept while the others are discarded. If “median” (resp. “mean”), then the threshold value is the median (resp. the mean) of the feature importances. A scaling factor (e.g., “1.25\*mean”) may also be used. If None and if available, the object attribute threshold is used. Otherwise, “mean” is used by default.

---

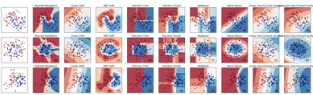
**Returns:**     **X\_r** : array of shape [n\_samples, n\_selected\_features]

The input samples with only the selected features.

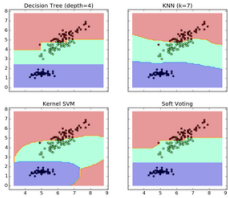
---

## Examples using sklearn.tree.DecisionTreeClassifier

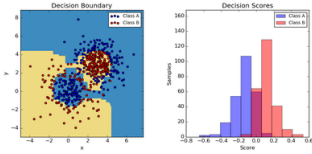




Classifier comparison

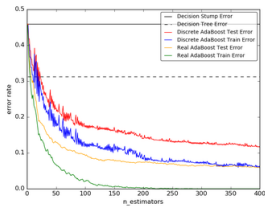


Plot the decision boundaries of a VotingClassifier

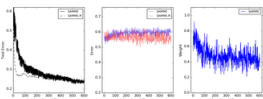


Two-class AdaBoost

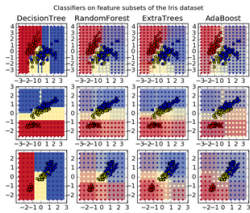
>>



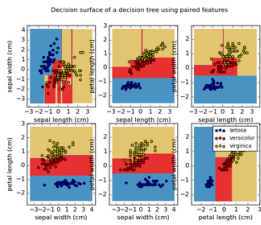
Discrete versus Real AdaBoost



Multi-class AdaBoosted Decision Trees



Plot the decision surfaces of ensembles of trees on the iris dataset



Plot the decision surface of a decision tree on the iris dataset